

SPIN: An Open Simulator of Realistic Spacecraft Navigation Imagery

Javier Montalvo

Video Processing and Understanding Lab, Universidad Autónoma de Madrid, Spain

Juan Ignacio Bravo Pérez-Villar

Deimos Space, Madrid, Spain. Video Processing and Understanding Lab, Universidad Autónoma de Madrid, Spain

Álvaro García-Martín

Video Processing and Understanding Lab, Universidad Autónoma de Madrid, Spain

Pablo Carballeira

Video Processing and Understanding Lab, Universidad Autónoma de Madrid, Spain

Jesús Bescós

Video Processing and Understanding Lab, Universidad Autónoma de Madrid, Spain

Abstract— Data acquired in space operational conditions is scarce due to the costs and complexity of space operations. This poses a challenge to learning-based visual-based navigation algorithms employed in autonomous spacecraft navigation. Existing datasets, which largely depend on computer-simulated data, have partially filled this gap. However, the image generation tools they use are proprietary, which limits the evaluation of methods to unseen scenarios. Furthermore, these datasets provide limited ground-truth data, primarily focusing on the spacecraft’s translation and rotation relative to the camera. To address these limitations, we present SPIN (SPacecraft Imagery for Navigation), an open-source realistic spacecraft image generation tool for relative navigation between two spacecrafts. SPIN provides a wide variety of ground-truth data and allows researchers to employ custom 3D models of satellites, define specific camera-relative poses, and adjust various settings such as camera parameters and environmental illumination conditions. For the task of spacecraft pose estimation, we compare the results of training with a SPIN-generated dataset against existing synthetic datasets. We show a %50 average error reduction in common testbed data (that simulates realistic space conditions). Both the SPIN tool (and source code) and our enhanced version of the synthetic datasets will be publicly released upon paper acceptance on GitHub. <https://github.com/vpulab/SPIN>.

Index Terms— Aerospace Navigation, Pose estimation, Simulation, Synthetic Data,

This work has been supported by the Ministerio de Ciencia, Innovación y Universidades of the Spanish Government under HVD (PID2021-125051OB-I00) and SEGA-CV (TED2021-131643A-I00) projects, and the Comunidad Autónoma de Madrid under the Grant IND2020/TIC-17515 *Corresponding Author J. Montalvo. J. Montalvo and JI. Bravo Pérez-Villar are co-first authors.*

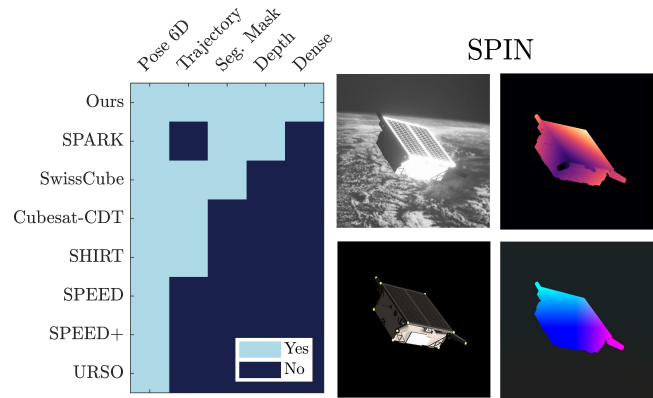


Fig. 1. A comparative between SPIN and other satellite datasets. On the left, the matrix compares their features (light blue cells indicate available features). Images show SPIN data examples: a rendering and a depth ground-truth (top images); our keypoint tool and a dense-pose ground-truth (bottom images).

I. Introduction

ACQUIRING data from space presents certain difficulties which include, but are not limited to, the high costs of designing and launching a spacecraft, the limited on-board resources for data storage and transmission, and the required precise spacecraft control to obtain the desired representative data. This limits the availability of data acquired in real-space operational conditions. While substantial efforts have been made to overcome these limitations for observation missions of the Earth [1], Sun [2], and other celestial bodies with high scientific value [3], [4], [5], there is still a noticeable gap in other space-related fields. Specifically, data to support autonomous spacecraft relative navigation remains scarce [6].

A particular case of interest in spacecraft relative navigation is the interaction between two non-cooperative spacecrafts. This is crucial for supporting current and future space missions, including tasks such as on-orbit servicing, active debris removal, close formation flying, rendezvous and docking, or space exploration. In all these mission scenarios, autonomy is indispensable as both the signal delays and the limited bandwidth render remote spacecraft operation unfeasible.

Visual-based navigation plays a crucial role in achieving autonomy, as it provides information on the elements of the environment and the relative position and orientation of the spacecraft. Current state-of-the-art methods for visual-based navigation employ learning techniques. However, the aforementioned lack of data prevents training robust learning-based navigation algorithms to support the required autonomy. To address this issue, researchers often turn to computer-based simulators that generate synthetic datasets, as detailed in Section II. These datasets are essential for research in visual-based spacecraft navigation: they have enabled the development of algorithms capable of estimating spacecraft pose with centimetre-level accuracy [7] even in situations with large differences in scale [8].

The tools and models used for synthetic data generation are proprietary, which poses two main constraints. First, there is a difficulty in adding ground-truth data to the existing datasets which typically lack depth information, segmentation information or dense pose annotations, i.e. mapping each pixel to its corresponding location on the 3D model. Second, these proprietary tools restrict the possibility to test models in scenarios not covered by the datasets. For instance, the datasets generally do not include videos, defined here as sets of temporally related images, even though such image sequences are frequent in real rendezvous scenarios. Finally, assessing models under specific conditions like reflective surfaces or varied backgrounds, which are vital for creating more realistic images, is also challenging due to these limitations.

To bridge this gap, we present the SPacecraft Imagery for Navigation (SPIN) tool, the first open-source image generation tool designed specifically to create data for visual-based navigation between two spacecrafts. SPIN addresses the limitations of current state-of-the-art datasets, as illustrated in Figure 1, by enabling the generation of highly realistic and customizable spacecraft imagery in various poses or pose sequences. It also offers extensive ground truth data, encompassing depth, dense pose annotations, segmentation, and keypoints. SPIN allows to load any external spacecraft 3D model and comes pre-loaded with a model of the Tango spacecraft.

To validate SPIN, we show experiments results on the task of pose estimation, the common benchmark across all identified datasets (see Figure 1). To do so, we compare the results obtained from training with the SPIN-generated images against those obtained with two widely-used datasets: SPEED+ [7] and SHIRT [9] further introduced in Section II. SPEED+ and SHIRT include both synthetic and testbed images –realistic imagery captured in laboratory conditions. We compare by training with a SPIN replication of the synthetic datasets, with the same poses but increased realism and ground-truth data, and then evaluating on the testbed images, yielding a 53% average reduction in error rate for spacecraft pose estimation tasks. We summarise our **contributions** as follows:

- We provide the first open-source simulation tool designed to generate realistic datasets of spacecraft images along with depth, segmentation and dense pose ground-truth data.
- Our tool allows to augment the ground-truth data of existing datasets, provided that the spacecraft’s 3D model is publicly available.
- We provide an *enhanced* version of the existing SPEED+ [7] and SHIRT [9] datasets, with improved realism and additional depth, segmentation and dense pose labels.

II. Related Work

In modern vision-based algorithms that employ Convolutional Neural Networks or Transformer architectures, data serves a pivotal role in facilitating effective training and achieving optimal performance. In the space operations domain, accumulating large datasets acquired in operational conditions is impractical, due to factors such as high costs, restricted on-board resources, or constrained communication links. Two primary approaches are employed to replace real space imagery: testbed facilities and computer-based simulators. Testbed facilities are specialised lab setups designed to mimic real conditions. In the context of relative spacecraft navigation, these facilities typically feature a scaled mock-up of the target spacecraft, a motion system (e.g., a robotic arm) to simulate spacecraft dynamics, authentic camera engineering models for imaging, specialised illumination systems, and the required control and computational infrastructure. With respect to simulators, rendering tools provide computer-generated images that emulate the visual characteristics of such navigation scenarios, also featuring adjustable camera parameters, customised lighting conditions, and tailored backgrounds. Additionally, they supply precise ground-truth data for aspects such as pose, depth, segmentation, and object detection.

There exists a trade-off among cost, flexibility, and representativeness when choosing between testbed and simulator approaches. Rendering tools offer a cost-effective and flexible solution, enabling the easy generation of diverse scenarios and backgrounds. In contrast, testbed facilities utilise real hardware and accurate mock-ups, yielding images that more closely resemble actual space-operational conditions, thereby reducing the domain gap. However, the substantial costs and specialised hardware requirements associated with testbed facilities restrict their widespread adoption in open research. Consequently, they are often reserved for secondary adaptation stages or for validation and verification purposes.

We provide a detailed description of publicly available datasets based on monocular intensity images, summarising their features and limitations in Table I. While our focus is on optical datasets, we acknowledge the existence of datasets derived from event sensors in the literature [10].

A. Datasets

The SPARK dataset [11] includes over 150,000 synthetic RGB images, along with corresponding depth and segmentation masks. It features 10 different spacecraft models obtained from NASA’s 3D resources and 5 distinct debris objects. The images are rendered with the Unity framework.

The SwissCube Dataset [8] constitutes a synthetic collection featuring 50,000 images of a 1U CubeSat model based on the SwissCube satellite. These images are organized into 500 trajectories, each comprising 100

TABLE I

Comparison of various established satellite pose estimation datasets with our simulation tool, SPIN. A checkmark (✓) indicates the availability of the feature, while a dash (-) indicates its absence. Pose 6D refers to the pose encoded with the absolute position and rotation, whereas dense coordinates indicate the 3D position of each pixel of the spacecraft. The depth column represents the ground-truth metric depth, while the segmentation column represents the segmentation mask. The Images column indicates the availability of gray or colour images and their resolution. The Simulation column indicates the range of distances between spacecraft and camera covered by the dataset, and the availability of trajectories (video sequence).

	Labels				Images		Simulation		
	Pose 6D	Dense Coord.	Depth	Segmentation	Bands	Resolution	Range	Trajectories	Testbed
SPARK [11]	✓	-	✓	✓	RGB	1024x1024	[1.5m, 10m]	-	-
SwissCube [8]	✓	-	-	✓	RGB	1024x1024	[0.1m, 1m]	✓	-
CubeSat-CDT [12]	✓	-	-	-	RGB	1440x1080	[0.4m, 3.8m]	✓	✓
URSO [13]	✓	-	-	-	RGB	1080x960	[10m, 40m]	-	-
SPEED [14]	✓	-	-	-	Gray	1920x1200	[3m, 40.5m]	-	✓
SPEED+ [7]	✓	-	-	-	Gray	1920x1200	≤ 10m	-	✓
SHIRT [9]	✓	-	-	-	Gray	1920x1200	≤ 8m	✓	✓
SPIN (Ours)	✓	✓	✓	✓	RGB/Gray	Configurable	Configurable	✓	-

frames, and are generated using the Mitsuba Renderer 2 framework. Importantly, the dataset encompasses a broad range of distances to the CubeSat and the camera, thereby introducing significant scale variability.

The CubeSat Cross-Domain Trajectory (CDT) dataset [12] comprises RGB images captured across multiple trajectories of a 1U CubeSat within three distinct domains. Specifically, the dataset includes two synthetic domains: the first, generated using Unity, contains 50 trajectories, and the second, created with Blender, encompasses 15 trajectories. Additionally, a testbed domain is provided, featuring 21 trajectories. The data from all three sources employ the same 1U CubeSat and share identical camera intrinsic parameters.

The Unreal Rendered Spacecraft On-Orbit (URSO) Dataset [13] is generated using Unreal Engine 4 and features a total of 15,000 synthetic RGB images. The dataset is divided into three distinct subsets, each containing 5,000 images. One subset focuses on the Dragon spacecraft, while the remaining two subsets present varying levels of complexity for the Soyuz spacecraft.

The Spacecraft Pose Estimation Dataset (SPEED) [14] contains 15,000 synthetic grayscale images, in addition to 305 images captured under testbed conditions. The dataset focuses on the Tango spacecraft from the PRISMA mission. Each image is annotated for pose estimation tasks.

The Next Generation Spacecraft Pose Estimation Dataset (SPEED+) [7], represented in the bottom row of Figure 2, consists of 60,000 synthetic images featuring the Tango spacecraft, accompanied by pose annotations. In addition, the dataset includes 9,531 annotated testbed images of a half-scale mock-up model. These test images are divided into two subsets: Sunlamp, which contains 2,791 images characterised by strong illumination and reflections against a dark background; and Lightbox, featuring 6,740 images with softer lighting conditions,

elevated noise levels, and the presence of the Earth in the background.

The Satellite Hardware-In-the-loop Rendezvous Trajectories Dataset (SHIRT) [9], represented in the top row of Figure 2, features two distinct trajectories (video sequences), ROE1 and ROE2, capturing the poses of a Tango satellite from the perspective of a service spacecraft. Each sequence offers two sets of images: synthetic grayscale images and hardware-in-the-loop testbed images that are similar to the Lightbox subset in the SPEED+ dataset. Both the synthetic and testbed images capture the spacecraft in identical poses and are acquired using the same camera intrinsic parameters. This allows for the direct evaluation of the domain gap impact while holding all other variables constant.

We employ SPEED+ and SHIRT as benchmarks due to their relevance in the spacecraft pose estimation literature. SPEED+ was employed in the European Space Agency Spacecraft Pose Estimation Challenge 2021 [7]. SHIRT expands SPEED+ to contain sequences of images.

B. Discussion

We argue that the current research in relative navigation between spacecrafts is constrained by existing datasets and the lack of simulation tools (Table I). Firstly, datasets generally lack of diverse ground-truth labels, not providing depth, segmentation, or dense pose information. For instance, dense prediction techniques for spacecraft pose estimation have demonstrated their efficacy [15], while only one dataset provides the means to compute dense pose via its dense depth maps [11]. Additionally, self-supervised methods for estimating monocular depth and pose [16], are hard to assess due to a lack of datasets with trajectories (video sequences) containing ground-truth depth.

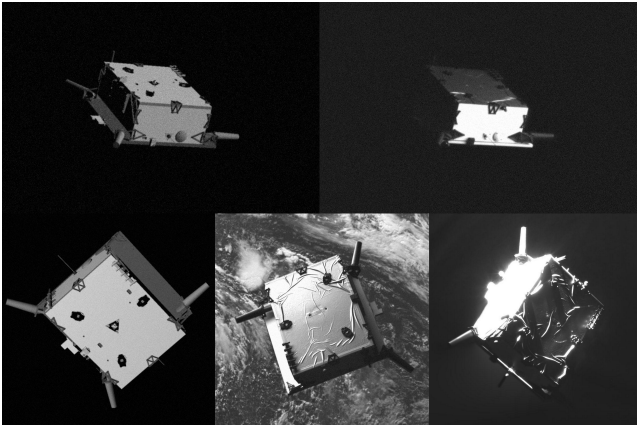


Fig. 2. Representative images from the SHIRT [9] and SPEED+ [7] datasets. In the top row images from the SHIRT dataset: the image on the left is from the synthetic domain, while the one on the right is from the testbed domain. The bottom row contains images from the SPEED+ dataset, are arranged from left to right, showcasing the simulated and the testbed –Lightbox, and Sunlamp– settings.

Secondly, the absence of image generation tools restricts the range of scenarios in which we can evaluate current state-of-the-art algorithms just to those considered by the existing datasets. Current datasets show limited variability in terms of changing surface reflectivity, camera modelling, background conditions, or variable sequences of poses. These limitations, for instance, prevent exploring research topics such as evaluating the effects of varying camera intrinsics –which could result from miscalibration during launch– on tasks such as monocular depth estimation [17] or pose estimation [18], [19].

III. SPIN: SPacecraft Imagery for Navigation

Figure 3 provides a schematic overview of our proposed tool, highlighting its essential features. The tool requires two primary inputs: a 3D model of a spacecraft (the Tango spacecraft is provided by default) and a set of predefined poses. Each pose defines the orientation and position of the spacecraft with respect to the camera so that the output is a set of images of the spacecraft in each defined pose with the associated ground-truth data.

Users have the flexibility to customise the rendered images acting on three main scene elements: 1) camera adjustments, specifically modifying intrinsic parameters, and sensor and lens imperfections; 2) altering the environment by tweaking illumination sources, background, and the rendering of shadows; and 3) activating or deactivating specific materials and material properties such as specularity.

The primary output of the tool are the spacecraft RGB or grayscale images, and the associated ground-truth data includes depth maps, dense pose maps, and segmentation masks. Additionally, a keypoint labelling tool to select keypoints and generate keypoint heatmaps is provided.

The tool has been developed in Unity Engine and will be distributed completely open-source.

A. Input

SPIN takes a 3D spacecraft model and a set of poses as inputs. We include a default 3D model of the Tango spacecraft, the one depicted in Figure 3. This model was created using the Fusion360 modelling software. There are multiple options for configuring spacecraft poses: manually setting the spacecraft’s location and rotation, importing camera-relative poses from a file, or generating random poses within the tool. In this last case, the tool employs the algorithm described in [20] to ensure uniform sampling of the rotation space.

B. Camera

The camera module governs image acquisition options and offers two configuration layers. The first layer focuses on lens distortion parameters and on the camera intrinsics, which determine parameters such as image size, focal length, and principal point. The second layer models specific challenges associated with space imaging, such as the absence of an atmosphere, low signal-to-noise ratio (SNR), expansive dynamic range, and highly reflective surfaces. These factors influence both the lens and the sensor. Within this layer, our tool provides options to simulate: a) sensor noise, b) lens glare and bloom coming from the sun and highly reflective surfaces, and c) different colour adjustment options to modify the overall colour temperature, intensity, and saturation of the final image.

C. Environment and Materials

The environment settings in our tool offer control over three main elements: illumination, background, and shadow rendering. For illumination, the tool includes three preset options: *Spotlight*, which simulates a spotlight in a fixed location and rotation with respect to the camera; *Sunlight*, offering directional lighting that replicates solar rays; and *Ambient Light*, producing diffuse illumination by combining multiple light sources around the target. All modes feature user-adjustable light intensity, the *Spotlight* and *Sunlight* options additionally allow to adjust rotation and position, and multiple Spotlight lights can be used and adjusted at the same time. Moreover, different illumination presets can also be enabled simultaneously.

As for the background, our tool not only supports a uniform backdrop but also includes a default selection of Earth images taken from space. These images can either be fixed or randomly positioned during each generation. Users have the flexibility to incorporate their own background images as well. Additionally, the tool provides the capability to toggle shadows and ambient occlusion on or off and allows for the customisation of various shadow properties.

Finally, the tool includes an option to enable or disable high-quality materials that offer enhanced textures and reflections. Given that spacecraft are often covered with

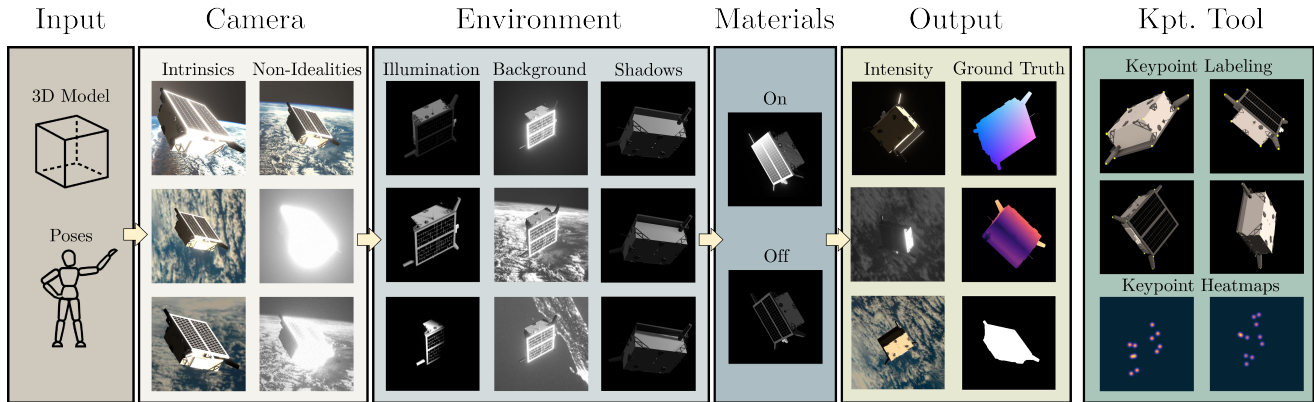


Fig. 3. Pipeline of SPIN. The input to the tool is a 3D model of the spacecraft (the Tango one is the default) and a set of poses. Realism can be configured acting on three scene elements: 1) Camera, that allows to modify the intrinsics and non-idealities such as camera glare, sensor noise, and color adjustment; 2) Environment, that allows to define illumination, background, and shadows rendering; and 3) Materials, that allow enabling high-quality reflective materials. SPIN outputs the intensity images (RGB or grayscale) and ground-truth data including the dense pose, the depth, and the segmentation mask; additionally, a keypoint labeling and heatmap generation tool is provided.

reflective materials, this feature allows to choose non-Lambertian surfaces, which ensure consistent illumination across various spacecraft poses.

In Figure 4 we include some examples of images that can be generated using our simulator, with different settings to the ones used for the SPEED+ like images.

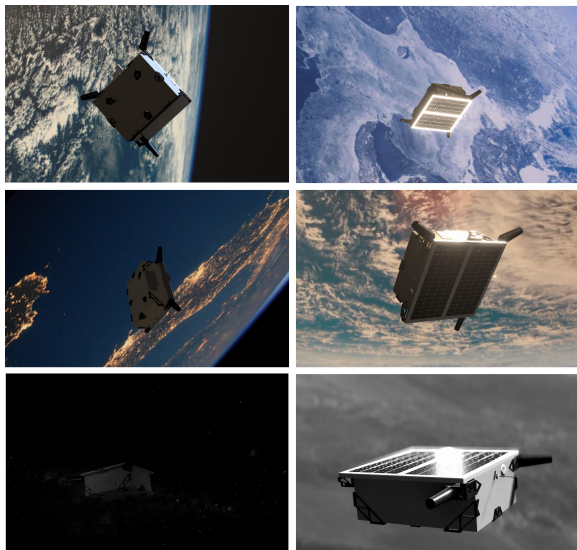


Fig. 4. Some examples of different renderings generated with our tool. Top right image depicts the keypoint visualization tool, whereas the rest of the images present different combinations of background, illumination, and color corrections.

D. Outputs

The main output is the set of intensity images based on the configured camera, environment, and materials. Additionally, the tool can provide a metric depth map of the scene, a segmentation mask, and dense 3D coordinates of the spacecraft. While their coordinates can be derived from the depth map, we provide them separately to simplify the workflow and to ease research with prediction

methods based on dense pose. In addition, we provide a way to generate keypoint heatmap for keypoint-based pose estimation tasks. This is detailed in Section E.

E. Keypoint Tool

The tool includes a keypoint labelling module and a keypoint heatmap generation module. Keypoint-based methods currently represent the state-of-the-art in spacecraft pose estimation [6]. These methods rely on estimating the 2D projection of predefined 3D keypoints (as heatmaps) on the spacecraft and solving for 2D-3D correspondences to determine the pose. Despite their effectiveness, few datasets provide a set of keypoints. The typical process involves selecting a set of 2D points, triangulating them, and further optimising the triangulation using convex solvers [21]. As we are aware that this is a time-consuming task, to facilitate research we provide a labelling mechanism for keypoints and a heatmap generation module to create the heatmap representation of the 2D keypoints.

IV. Experimental Validation

We propose to validate SPIN by employing its output to train a model for the task of spacecraft pose estimation, the common benchmark across the existing datasets. More in detail, we use SPIN to replicate the synthetic images of SPEED+ [7] and SHIRT [9] for our validation (examples provided in Figure 2). SPEED+ shows the particularity of having a synthetic and two testbed domains which, in addition to the testbed domain of SHIRT, provide a suitable framework for evaluating the quality of SPIN and its impact on sim-to-real transfer. The experiments are organised as follows, we first describe the settings for the experiments on the pose estimation. Next, we detail the configurations used in SPIN for image generation. We conclude by presenting the quantitative results, including an ablation study.

TABLE II

Summary of available parameters and options in the simulator. The two right-most columns indicate which parameters have been modified for increased realism and which ones have been set to preserve a consistent input domain w.r.t the original synthetic split of SPEED+ and SHIRT datasets.

Module	Parameter	Description	Options	Modification w.r.t SPEED+ SHIRT	Justification
Camera	Intrinsics + Lens distortion	Camera sensor and lens configuration	Sensor Size, FOV, ISO, Aperture...	No	Preserve consistent input domain
	Color Adjustment	Enables different color adjustment options over the generated image	Saturation, Contrast, Hue Shift...	No	Preserve consistent input domain
	Noise	Simulates different types of sensor noise	Noise type, intensity, response.	No	Preserve consistent input domain
	Glare/Bloom	Simulates camera glare effects.	Threshold, Intensity, Scatter, Tint...	Yes	Enhance scene realism to align with testbed conditions
Environment	Illumination	Multiple illumination options.	Ambiental Illumination, Directional, Spot-light, Light intensity, Light temperature...	Yes	Enhance scene realism to align with testbed conditions
	Background	Allows the user to set up background images.	Background image, position, rotation, random positions.	No	Maintain a consistent Earth background across respective poses for uniform impact
	Shadows	Enables shadows on the simulator.	Shadow quality, Shadow distance, Ambient Occlusion...	No	Retain original shadow rendering
Material	High-quality materials	Enable or disable high-quality materials for the satellite.	Material quality, Specularity on/off.	Yes	Enhance scene realism to align with testbed conditions

A. Pose Estimation Experimental Settings

In all our experiments we consistently use the same architecture and the same evaluation metrics, described in Section 1 and Section 2 respectively. All the models described are trained using PyTorch [22] with input 512x512 grayscale images. We employ the Adam optimiser [23] with a learning rate of 0.0001. The ground-truth heatmaps are created with SPIN always using the same parameters (sigma deviation of 7 pixels) and all training parameters are kept the same for all models, including the random seed. This approach is adopted to reduce the influence of factors other than the input training data.

In the SPEED+ experiments, we train two distinct models: the first on the original synthetic SPEED+ training split, which comprises 47,966 images. The second on the corresponding images replicated using the SPIN tool. Both models are trained for 60 epochs and are tested over the testbed domains of SPEED+: Sunlamp with 2,791 images, and Lightbox with 6,740.

For the SHIRT experiments, we train four separate models. Two are trained on the original synthetic sequences ROE1 and ROE2, each containing 2,371 images. The other two models are trained on replicas of these sequences, created using the SPIN tool. These four models are trained for 5 epochs and are evaluated in the testbed domains of ROE1 and ROE2. These testbed domains consist of the same poses as those in their respective training sequences, but the testing is conducted in a different environment specific to the testbed domains.

1. Pose Estimation Model

We choose a simple baseline model from [24] to capture the performance differences introduced by the different training data. Given a spacecraft image, we use a ResNet-50-based architecture to regress a heatmap $\hat{h} \in \mathbb{R}^{N \times M \times C}$, where N and M are the image dimensions and C is the number of unique keypoints. Each channel c of \hat{h} encodes a 2D Gaussian heatmap centered at the predicted image 2D coordinates \hat{p}_i corresponding

to each spacecraft 3D keypoint P_i . The ground-truth keypoint positions p_i to generate the ground-truth heatmap h are computed by projecting P_i using the ground-truth pose T with the perspective equation. The network is trained to minimise the mean squared error between \hat{h} and h , as given by:

$$\ell_h = \frac{1}{NMC} \sum \|\hat{h} - h\|_F^2. \quad (1)$$

At test time, the estimated keypoint coordinates \hat{p}_i are determined by locating the maximum value in the i^{th} channel of \hat{h} . Finally, we employ an EPnP method [25] within a RANSAC loop to retrieve the pose estimate, using the 2D-3D correspondences and the camera intrinsic parameters.

2. Metrics

We adopt the evaluation metrics defined in [26]. The translation error E_v is calculated as the Euclidean distance between the estimated translation vector \hat{v} and its ground-truth counterpart v , formulated as $E_v = \|\hat{v} - v\|_2$. Similarly, the orientation error E_q is determined by the rotation angle required to align the estimated quaternion \hat{q} with the ground-truth quaternion q , given by $E_q = 2 \cdot \arccos(|\langle \hat{q}, q \rangle|)$. These errors are subsequently converted into scores: the translation score is $S_v = E_v / \|v\|_2$, and the orientation score is $S_q = E_q$. Any translation and orientation scores falling below 2.173×10^{-3} and 0.169° , respectively, are set to zero [26]. The total score is then computed as $S = S_q + S_v$.

B. SPIN Generation Settings

For the generation of the dataset replicas with SPIN we set some parameters to match those of SPEED+ and SHIRT with the aim of keeping a consistent input domain, and we set others to improve the image realism. We provide a summary of the modified parameters in Table II. Specifically, for camera settings, we keep the intrinsics, color adjustment, and noise parameters constant. This

TABLE III

Comparison of spacecraft pose estimation performance between the original SPEED+ dataset and the SPEED+ replica generated with SPIN. Lower scores indicate better performance.

Training data	Test	S_v	S_q	S
SPEED+	Lightbox	0.359	1.361	1.715
	Sunlamp	0.356	1.665	2.021
SPIN	Lightbox	0.171	0.719	0.891
	Sunlamp	0.122	0.621	0.743

TABLE IV

Comparison of spacecraft pose estimation performance between the original SHIRT dataset over the ROE1 and ROE2 sequences and the respective replicas generated with SPIN.

	ROE	S_v	S_q	S
SHIRT	1	0.072	0.542	0.614
	2	0.526	2.131	2.021
SPIN	1	0.030	0.352	0.381
	2	0.176	0.987	1.163

ensures that the spacecraft is viewed from the same perspective and that noise levels are uniform, for fair comparison. However, we modify glare and bloom settings to enhance the scene’s realism. For the environment settings, we use shadow rendering techniques similar to those used in the SPEED+ dataset. To minimize possible effects of Earth’s presence in the background for the task of pose estimation, we choose to use similar background images in our synthetic dataset as those in SPEED+. Regarding illumination, we adjust the settings to produce more realistic images, in particular, we employ the *Sunlight* setting. Lastly, we activate the high-quality materials to create images that more closely resemble those from the testbed domain. In Figure 7 we display examples of the differences between images generated using a configuration aligned with the SPEED+ synthetic subset, and the same images after applying our enhancement settings.

C. Quantitative Results

Table III indicates that algorithms trained on our synthetic dataset exhibit significant improvements over those trained on the SPEED+ dataset: a 48% reduction in the error rates for the Lightbox testbed, moving from a baseline error of 1.71 to a 0.891; and a 63% drop for the Sunlamp one, from 2.021 to 0.743. For the SHIRT dataset, results in Table IV indicate a 47% error rate reduction for ROE1, and a 53% one for ROE2. Visual examples of the pose estimation errors are presented in Figure 5.

These experiments show how SPIN allows to reduce the performance gap between the synthetic and testbed

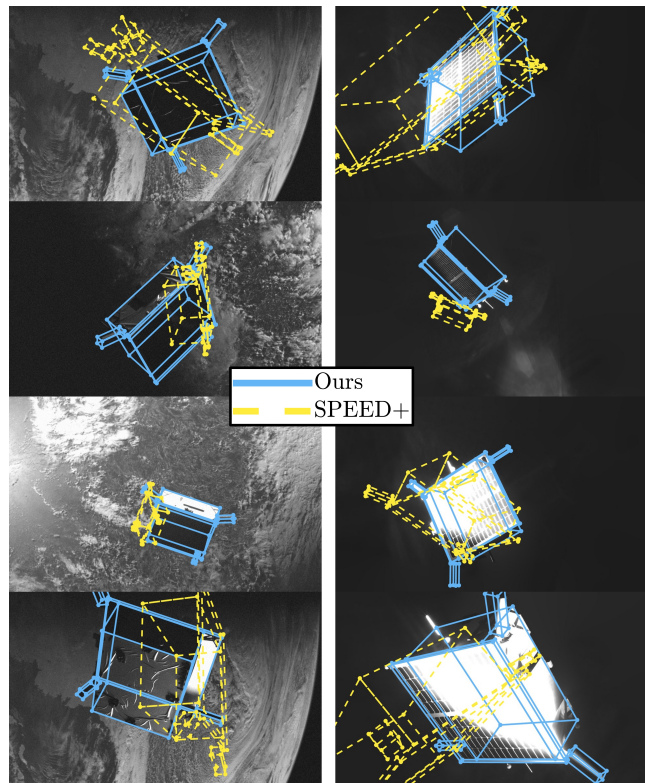


Fig. 5. Examples of the estimated pose represented as a wireframe model (solid blue for SPIN and dashed yellow for SPEED+) overlaid over the real satellite for SPEED+ images of the Lightbox (left) and Sunlamp (right) testbeds.

domains, by providing more realistic training images than those of the existing synthetic datasets from the literature.

1. Ablation Tests

In this section, we include an ablation study to evaluate the specific impact of the Camera, Environment, and Material settings modifications for improved realism on the pose estimation performance. We conduct the study by replicating different synthetic versions of the training split of the SPEED+ with SPIN. Each replica is created with a different combination of SPIN settings and used to train a separate pose estimation model. We illustrate the effect of the SPIN settings on a sample image in Figure 6. As in previous experiments, we evaluate over the testbed SPEED+ domains of Sunlamp and Lightbox.

Table V summarises the results. We first define a baseline –first row– by deactivating all settings. Next, we evaluate each setting independently. Activating each setting independently does not result in significant performance improvements, probably due to the lack of an intense directional light to trigger camera effects such as lens bloom or better reflective materials, hence generating images very similar to the baseline ones. The activation of the Environment settings leads to the most effective results, allowing for the generation of more *challenging* images, as the shadows cast by the directional light occlude parts of the satellite.

TABLE V

Impact of SPIN settings (indicated by a checkmark (✓) when activated and a dash (-) when deactivated) on the pose estimation performance over the SPEED+ testbeds. We activate or deactivate only those parameters that are configured differently from SPEED+, as indicated in Table II.

Settings			Lightbox			Sunlamp		
Camera	Environment	Materials	S_v	S_q	S	S_v	S_q	S
-	-	-	0.291	1.260	1.551	0.379	1.808	2.187
✓	-	-	0.352	1.450	1.803	0.416	1.793	2.209
-	✓	-	0.309	1.269	1.578	0.224	1.146	1.370
-	-	✓	0.381	1.483	1.864	0.437	1.775	2.213
✓	✓	-	0.217	0.911	1.13	0.132	0.759	0.891
-	✓	✓	0.353	1.508	1.861	0.392	1.657	2.049
✓	-	✓	0.270	1.093	1.364	0.215	1.103	1.319
✓	✓	✓	0.171	0.719	0.891	0.122	0.621	0.743

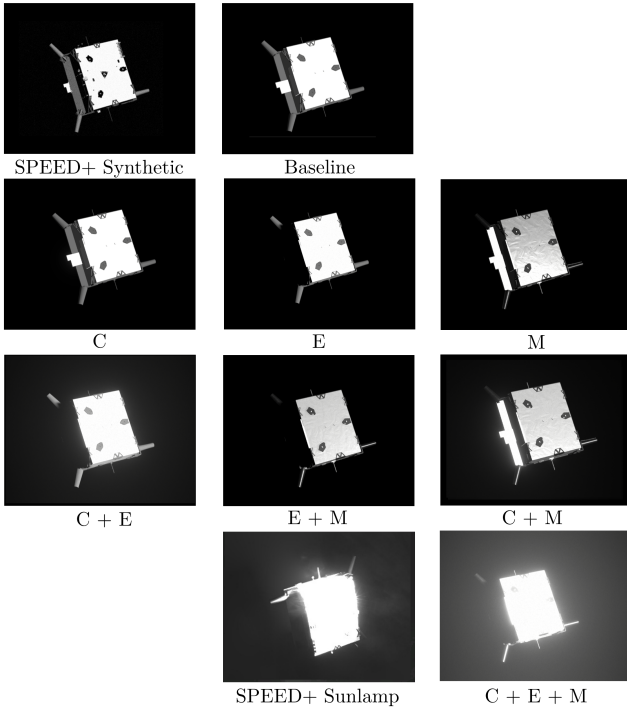


Fig. 6. Comparison of a synthetic SPEED+ image and those generated with SPIN under different configurations: *Baseline* indicates that no setting is activated, and “C,” “E,” and “M” indicate that settings for the Camera, Environment, and Material are enabled, respectively.

Combining couples of settings yields better results. Interestingly, when Materials and Environment settings are activated, while results improve with respect to just activating Materials, they do not exceed those of just activating Environment. Best results are clearly obtained by combining Camera and Environment. This synergy might be attributed to the increased illumination on the spacecraft due to either increased reflectivity or to a more

intense light source, which coupled with camera settings such as bloom, results in much more realistic images. It is noteworthy, however, that this same reflectivity and intense illumination can diminish performance when these camera settings are not enabled, leading to overexposed images.

Finally, enabling all three settings results in the best performance, reaching the results presented in Section C.

V. Conclusions

In conclusion, SPIN significantly enhances realism in image generation, narrowing the gap between synthetic and real imagery in pose estimation compared to existing synthetic datasets. It also provides a wider range of ground-truth data, including the use of dense pose labels, that no other dataset in the literature provides. Moreover, SPIN facilitates the creation of new and diverse test scenarios, expanding the variety and depth of ground-truth data in current datasets.

REFERENCES

- [1] P.-P. Mathieu and C. Aurbrecht, *Earth observation open science and innovation*. Springer Nature, 2018.
- [2] D. Müller, O. S. Cyr, I. Zouganelis, H. R. Gilbert, R. Marsden, T. Nieves-Chinchilla, E. Antonucci, F. Auchère, D. Berghmans, T. Horbury *et al.*, “The solar orbiter mission-science overview,” *Astronomy & Astrophysics*, vol. 642, p. A1, 2020.
- [3] D. L. Matson, L. J. Spilker, and J.-P. Lebreton, “The cassini/huygens mission to the saturnian system,” *Space Science Reviews*, vol. 104, no. 1-4, pp. 1–58, 2002.
- [4] S. J. Bolton, J. Lunine, D. Stevenson, J. Connerney, S. Levin, T. Owen, F. Bagenal, D. Gautier, A. Ingersoll, G. Orton *et al.*, “The juno mission,” *Space Science Reviews*, vol. 213, pp. 5–37, 2017.
- [5] S. Stern, H. Weaver, J. Spencer, H. Elliott, and N. H. Team, “The new horizons kuiper belt extended mission,” *Space Science Reviews*, vol. 214, pp. 1–23, 2018.

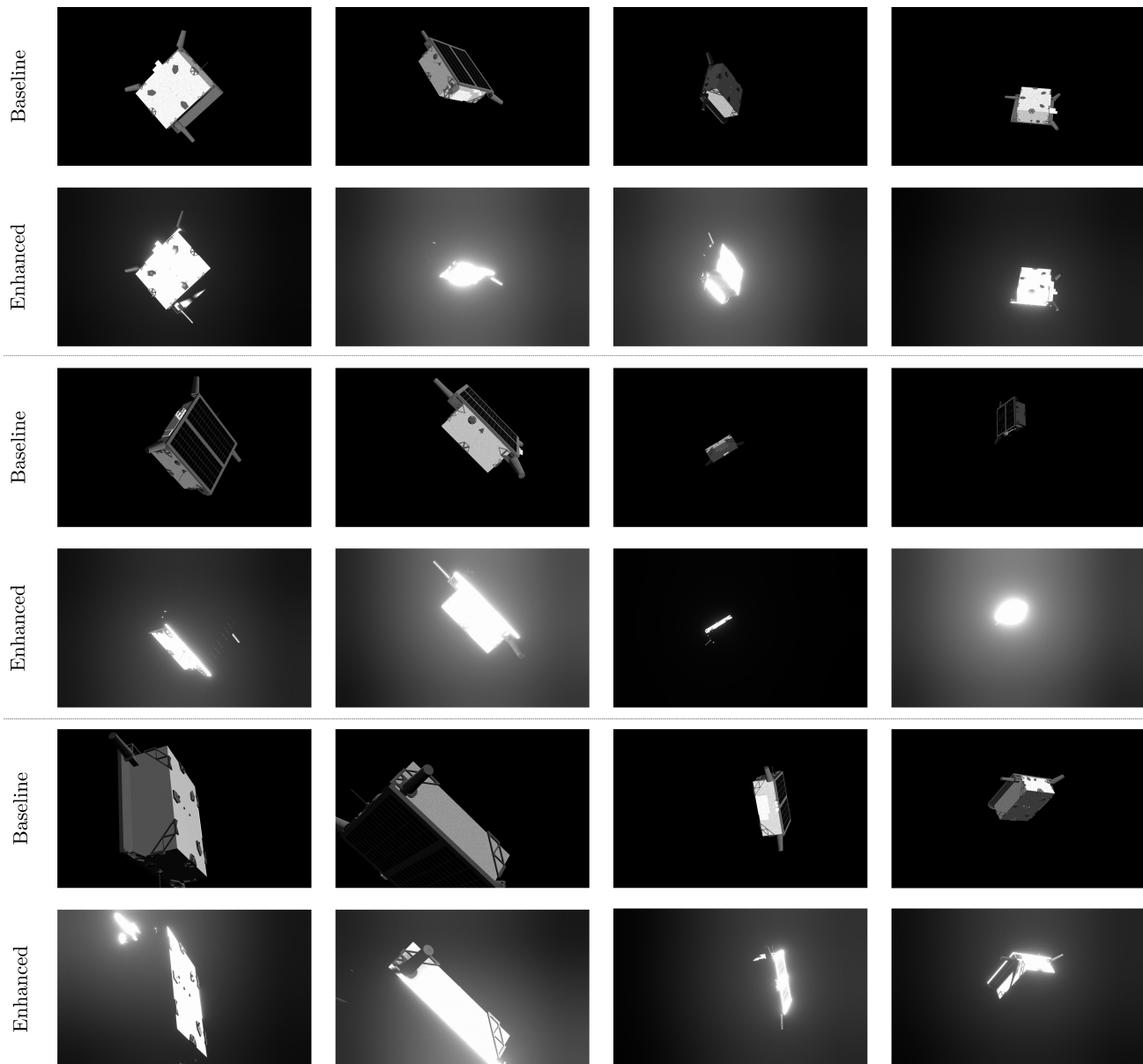


Fig. 7. Image comparison between a baseline configured like the synthetic split from the SPEED+ dataset, and the same images after applying Environment, Camera, and Material settings.

- [6] J. Song, D. Rondao, and N. Aouf, "Deep learning-based spacecraft relative navigation methods: A survey," *Acta Astronautica*, vol. 191, pp. 22–40, 2022.
- [7] T. H. Park, M. Märtens, G. Lecuyer, D. Izzo, and S. D'Amico, "Speed+: Next-generation dataset for spacecraft pose estimation across domain gap," in *2022 IEEE Aerospace Conference (AERO)*. IEEE, 2022, pp. 1–15.
- [8] Y. Hu, S. Speierer, W. Jakob, P. Fua, and M. Salzmann, "Wide-depth-range 6d object pose estimation in space," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 15 870–15 879.
- [9] T. H. Park and S. D'Amico, "Adaptive neural-network-based unscented kalman filter for robust pose tracking of noncooperative spacecraft," *Journal of Guidance, Control, and Dynamics*, vol. 46, no. 9, pp. 1671–1688, 2023.
- [10] M. Jawaid, E. Elms, Y. Latif, and T.-J. Chin, "Towards bridging the space domain gap for satellite pose estimation using event sensing," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023, pp. 11 866–11 873.
- [11] M. A. Musallam, V. Gaudilliere, E. Ghorbel, K. Al Ismaeil, M. D. Perez, M. Poucet, and D. Aouada, "Spacecraft recognition leveraging knowledge of space environment: simulator, dataset, competition design and analysis," in *2021 IEEE International Conference on Image Processing Challenges (ICIPC)*. IEEE, 2021, pp. 11–15.
- [12] M. A. Musallam, A. Rathinam, V. Gaudillière, M. O. d. Castillo, and D. Aouada, "Cubesat-cdt: A cross-domain dataset for 6-dof trajectory estimation of a symmetric spacecraft," in *European Conference on Computer Vision*. Springer, 2022, pp. 112–126.
- [13] P. F. Proença and Y. Gao, "Deep learning for spacecraft pose estimation from photorealistic rendering," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2020, pp. 6007–6013.
- [14] S. Sharma, T. Park, and S. D'Amico, "Spacecraft pose estimation dataset (speed)," *Stanford Digital Repository*, 2019.
- [15] M. Ulmer, M. Durner, M. Sundermeyer, M. Stoiber, and R. Triebel, "6d object pose estimation from approximate 3d models for orbital robotics," in *2023 IEEE/RSJ International Conference*

- on *Intelligent Robots and Systems (IROS)*. IEEE, 2023.
- [16] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1851–1858.
- [17] J. M. Facil, B. Ummerhofer, H. Zhou, L. Montesano, T. Brox, and J. Civera, “CAM-ConvS: Camera-Aware Multi-Scale Convolutions for Single-View Depth,” in *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [18] K. Josephson and M. Byrod, “Pose estimation with radial distortion and unknown focal length,” in *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2009, pp. 2419–2426.
- [19] V. Larsson, Z. Kukelova, and Y. Zheng, “Camera pose estimation with unknown principal point,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2984–2992.
- [20] J. J. Kuffner, “Effective sampling and distance metrics for 3d rigid body path planning,” in *IEEE International Conference on Robotics and Automation, 2004. Proceedings. ICRA’04. 2004*, vol. 4. IEEE, 2004, pp. 3993–3998.
- [21] B. Chen, J. Cao, A. Parra, and T.-J. Chin, “Satellite pose estimation with deep landmark regression and nonlinear pose refinement,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [22] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *Advances in neural information processing systems*, vol. 32, 2019.
- [23] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [24] B. Xiao, H. Wu, and Y. Wei, “Simple baselines for human pose estimation and tracking,” in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 466–481.
- [25] V. Lepetit, F. Moreno-Noguer, and P. Fua, “Epnnp: An accurate o (n) solution to the pnp problem,” *International journal of computer vision*, vol. 81, no. 2, pp. 155–166, 2009.
- [26] T. H. Park, M. Märtens, M. Jawaid, Z. Wang, B. Chen, T.-J. Chin, D. Izzo, and S. D’Amico, “Satellite pose estimation competition 2021: Results and analyses,” *Acta Astronautica*, vol. 204, pp. 640–665, 2023.



Javier Montalvo graduated with dual degrees in Computer Science and Mathematics, followed by a Master’s in Deep Learning for Audio and Video Processing from Universidad Autónoma de Madrid. He is currently a PhD student at the Video Processing and Understanding Lab (VPU-Lab) within the same institution. His research primarily focuses on the generation and utilization of synthetic data in machine learning applications.

The third paragraph begins with the author’s preferred title and last name (e.g., Dr. Smith, Prof. Jones, Mr. Kajor, Ms. Hunter, Mx. Riley). List any memberships in professional societies other than the IEEE. Finally, list any awards and work for IEEE committees and publications.



Juan Ignacio Bravo Pérez-Villar obtained a degree in Telecommunications Engineering in 2015 and received the titles belonging to the International Joint Master Program in Image Processing and Computer Vision (IPCVCV) in 2017 at the universities of Péter Pazmany (Hungary), Université de Bordeaux (France) and Universidad Autónoma de Madrid (Spain). Currently he is a PhD Student at the Video Processing and Understanding Lab (VPU-Lab)

and a Research Engineer of the GNC/AOCS Competence Centre at DEIMOS Space. His research interests are related to visual based navigation and domain adaptation.



Dr. Álvaro García-Martín received the M.S. degree in Electrical Engineering (“Ingeniero de Telecomunicación” degree) in 2007 (2002-2007) and the MPhil degree in Electrical Engineering and Computer Science (postgraduate Master) in 2009 and PhD in Computer Science in 2013 at Universidad Autónoma de Madrid (Spain). From 2006 to 2019, he has been with the Video Processing and Understanding Lab (VPU-Lab) at Universidad Autónoma of

Madrid as a researcher and teaching assistant. In 2008 he received a FPI research fellowship from Universidad Autónoma de Madrid. He is an Associate Professor (PhD) at Universidad Autónoma of Madrid since 2019. He has participated in several projects dealing with multimedia content transmission (PROMULTIDIS y MESH), video-surveillance (ATI@SHIVA) and activity recognition (SEMANTIC, EVENTVIDEO, HAVideo, HVD and SEGA-CV). He also serves as a reviewer for several international Journals (IEEE TIFS, IEEE CSVT, Springer MTAP,...) and Conferences (IEEE ICIP, IEEE AVSS,...). He has published more than 17 journal and conference papers. His current research interests are focused in the analysis of video sequences for the video surveillance (moving object extraction, object tracking and recognition, event detection...).



Dr. Pablo Carballeira received the Telecommunication Engineering degree (five years engineering program) in 2007 from the Universidad Politécnica de Madrid (UPM). He received the Communications Technologies and Systems Master degree (two year MS program) and the Ph.D. degree in Telecommunication from the UPM in 2010 and 2014 respectively. Since October 2017 he is an Assistant Professor in Universidad Autónoma de Madrid (UAM) and a member of the Video Processing and Understanding Lab (VPULab). His research interests include image processing and video coding, focusing on multiview and free-navigated video, computer vision and compressed sensing. He has been involved since 2008 in the standardization activities from ISO’s Moving Picture Experts Group (MPEG), related to multiview and free-navigated video.



Prof. Jesús Bescós received the Ingeniero de Telecomunicación degree in 1993 and a PhD in Communications in 2001, both at Universidad Politécnica de Madrid (UPM). He was member of the Image Processing Group at UPM (1993–2002), Assistant Lecturer at this University (1997–2002), and he is now Associate Professor at Universidad Autónoma de Madrid (since 2003), where he co-leads the Video Processing and Understanding Lab. His

research lines include video sequence analysis, content-based video indexing, 2D and 3D computer vision, etc. He has been actively involved in EU projects dealing with cultural heritage (e.g., RACE-Rama), education (e.g., ET-Trends), content analysis (e.g., ACTS-Hypermedia, ICTS-AceMedia and Mesh), virtual reality (e.g., IST-Slim-VRT), etc., leading to over forty publications in scientific conferences and journals. He is regular evaluator of national and international project proposals, and of submitted papers to conferences (ICIP, ICASSP, CBMI, SAMT, etc.) and journals (IEEE Tr. on CSVT, Tr. on Multimedia, Tr. on Image Processing, etc.)